

# FEDCVD: THE FIRST REAL-WORLD FEDERATED LEARNING BENCHMARK ON CARDIOVASCULAR DISEASE DATA

**Yukun Zhang<sup>1,3</sup>**   **Guanzhong Chen<sup>1</sup>**   **Zenglin Xu<sup>2,3</sup>**   **Jianyong Wang<sup>4</sup>**

**Dun Zeng<sup>5</sup>**   **Junfan Li<sup>1</sup>**   **Jinghua Wang<sup>1</sup>**   **Yuan Qi<sup>2,3</sup>**   **Irwin King<sup>6</sup>**

<sup>1</sup>Harbin Institute of Technology, Shenzhen, Shenzhen, China

<sup>2</sup>Fudan University, Shanghai, China

<sup>3</sup>Shanghai Academy of Artificial Intelligence for Science, Shanghai, China

<sup>4</sup>Sichuan University, Chengdu, China

<sup>5</sup>University of Electronic Science & Technology of China, Chengdu, China

<sup>6</sup>The Chinese University of Hong Kong, Hong Kong, China

yukun.zhang.cs@gmail.com   muxichenz@outlook.com   zenglin@gmail.com

wjy@scu.edu.cn   zengdun@std.uestc.edu.cn   {lijunfan, wangjinghua}@hit.edu.cn

qiyuan@fudan.edu.cn   king@cse.cuhk.edu.hk

## ABSTRACT

Cardiovascular diseases (CVDs) are currently the leading cause of death worldwide, highlighting the critical need for early diagnosis and treatment. Machine learning (ML) methods can help diagnose CVDs early, but their performance relies on access to substantial data with high quality. However, the sensitive nature of healthcare data often restricts individual clinical institutions from sharing data to train sufficiently generalized and unbiased ML models. Federated Learning (FL) is an emerging approach, which offers a promising solution by enabling collaborative model training across multiple participants without compromising the privacy of the individual data owners. However, to the best of our knowledge, there has been limited prior research applying FL to the cardiovascular disease domain. Moreover, existing FL benchmarks and datasets are typically simulated and may fall short of replicating the complexity of natural heterogeneity found in realistic datasets that challenges current FL algorithms. To address these gaps, this paper presents the first real-world FL benchmark for cardiovascular disease detection, named FedCVD. This benchmark comprises two major tasks: electrocardiogram (ECG) classification and echocardiogram (ECHO) segmentation, based on naturally scattered datasets constructed from the CVD data of seven institutions. Our extensive experiments on these datasets reveal that FL faces new challenges with real-world non-IID and long-tail data. The code and datasets of FedCVD are available <https://github.com/SMILELab-FL/FedCVD>.

## 1 Introduction

Cardiovascular Diseases (CVDs) cause over 18 million deaths globally each year, positioning them as one of the most significant global health challenges [1]. Early detection and accurate diagnosis of CVDs are crucial, as they allow for timely medical interventions and more effective treatment plans, which in turn significantly lower patient mortality rates [2]. In recent years, with the growing availability of electronic health records and other high-quality clinical data, researchers have increasingly utilized machine learning techniques to automate clinical diagnostics [3; 4], a strategy that has proven highly effective in the context of CVDs [5; 6]. This data-driven approach not only facilitates efficient early screening but also optimizes the allocation of healthcare resources, improving overall patient outcomes.

Compared to models trained in isolation at a single center, collaborations across multiple medical institutions enable the utilization of richer regional and demographic characteristics, fostering more precise and comprehensive research outcomes. However, medical data is considered highly sensitive, and recent privacy regulations (e.g., EU General Data Protection Regulation (GDPR) [7]) restrict its transfer, hindering the expansion of datasets through data sharing among institutions to train more efficient models, i.e., data isolation.

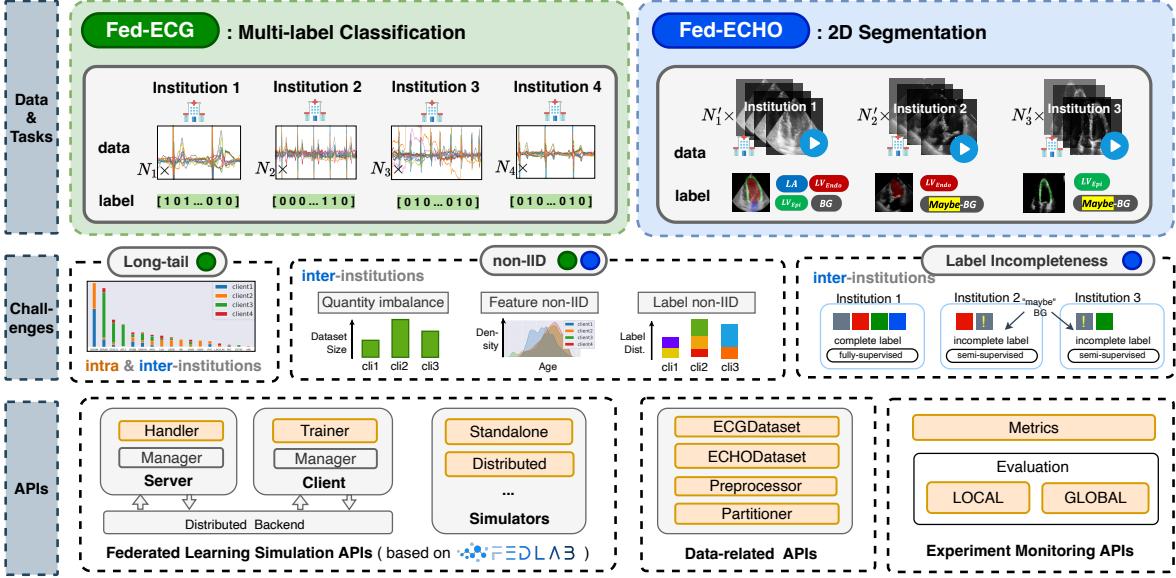


Figure 1: The overall architecture of the proposed FedCVD benchmark. We present two main settings (Fed-ECG, Fed-ECHO) and an experimental platform, highlighting three primary challenges. Green and blue circles in the challenges section indicate their presence in Fed-ECG and Fed-ECHO, respectively. The API section highlights user-facing APIs in orange boxes.

To address this issue, federated learning (FL) [8; 9] has been proposed as a more secure paradigm of distributed machine learning. A typical FL architecture involves a coordinator (Server) and several participants (Clients) with private data. By aggregating (e.g., FedAvg [9]) the model parameters or gradients from different Clients on the Server, participants collaboratively train high-performance models keeping private data within their respective domains. This process only involves transmitting model parameters, thus ensuring a certain degree of data privacy.

In medical applications such as CVD, integrating FL enables medical institutions to harness larger and more diverse datasets, collaboratively training models that are more unbiased and generalized, thereby enhancing diagnostic accuracy and clinical decision-making in real-world settings. For instance, [10] applied FL for joint case analysis across three institutions during the Covid-19 pandemic, significantly improving CT segmentation performance and facilitating more accurate detection of Covid-19. Additionally, the effectiveness of FL has been demonstrated in multi-center research, such as a study involving three centers focused on the medical image analysis task of whole prostate segmentation [11], further underscoring its relevance in realistic, large-scale medical scenarios.

Facilitating the application of FL in multi-center medical research, particularly in areas such as CVD, necessitates the creation of appropriate datasets and benchmarks to support the development of robust algorithms. However, publicly available cardiovascular disease datasets are limited, and those that do exist often suffer from incompatibility due to variations in data collection protocols. Furthermore, there is currently no comprehensive, publicly accessible benchmark specifically designed for FL on CVD data, which significantly impedes research progress in this domain. Additionally, most existing FL benchmarks simulate an FL scenario by manually partitioning data—often without considering geographic distribution—into smaller subsets, resulting in an overly idealized model that fails to capture the complexities and heterogeneity of real-world, multi-center CVD scenarios. This gap presents substantial challenges for the development and validation of effective FL algorithms in practical medical applications.

To address these gaps, we introduce the *first* multi-center FL benchmark specifically designed for CVD tasks, named **FedCVD**. Built from real-world CVD data collected from seven medical institutions (i.e., clients, the two terms will be used interchangeably), FedCVD utilizes a *natural partitioning* strategy. It comprises two primary datasets along with their corresponding tasks: electrocardiogram (ECG) classification and echocardiogram (ECHO) segmentation. FedCVD encapsulates three critical traits of FL in real-world CVD applications, each of which presents substantial challenges to FL algorithms:

**Challenging Trait 1. Non-IID Data:** The Non-independently and identically distributed (non-IID) data among institutions, including non-IID feature (e.g., variations in imaging quality due to different equipment across institutions)

Table 1: Comparison of FedCVD with Other Federated Datasets or Benchmarks.

	Long-tailedness Considered	Natural Partition	Incomplete Label	Covers CVD Domain	Code Available
FedDTI [12]	✗	✗	✗	✗	✓
FedTD [13]	✗	✗	✗	✗	✗
Flamby [14]	✗	✓	✗	✗	✓
NIPD [15]	✗	✓	✗	✗	✓
FEDLEGAL [16]	✓	✓	✗	✗	✓
FLHCD [17]	✗	✓	✗	✓	✓
FedMultimodal [18]	✗	*	✗	✗	✓
FedAudio [19]	✗	*	✗	✗	✓
<b>FedCVD</b>	✓	✓	✓	✓	✓

\*: Some datasets included are naturally partitioned.

and non-IID label (e.g., differences in disease prevalence across regions). The non-IID data may significantly hinder global model convergence.

**Challenging Trait 2. Long-tail Distribution:** The labels of CVD data from various institutions exhibit a long-tailed distribution, where a few labels dominate while most labels are sparse. This challenges the model’s performance on tail classes, a problem that is exacerbated in FL scenarios.

**Challenging Trait 3. Label Incompleteness:** For the same type of medical images, hospitals with strong annotation capabilities can identify all key segmentation areas, while those with weaker capabilities can identify only some. This incomplete annotation can mislead the global model’s segmentation performance in areas unrecognized by certain institutions.

Focusing on these challenging traits, FedCVD provides new insights and evaluation metrics for designing FL algorithms in multi-center CVD scenarios. Our contributions are summarized as follows:

1. We introduce FedCVD, an open-source federated multi-center healthcare dataset and benchmark specifically for the CVD domain. To the best of our knowledge, FedCVD is the largest multi-center CVD benchmark available. This dataset encompasses two critical tasks—multi-label classification and segmentation—within the CVD domain and includes data of varying scales. Crucially, all datasets are partitioned using natural splits.
2. Our benchmark emphasizes three critical traits in the FL-CVD scenario: non-IID, long tail, and label incompleteness. These traits pose significant challenges to existing FL algorithms.
3. We conducted extensive experiments on FedCVD to evaluate the performance of mainstream FL and centralized learning methods, validating the effectiveness of FL in the CVD context and the proposed three challenges. Additionally, we have made the open-source code in <https://github.com/SMILELab-FL/FedCVD> accessible for benchmark reproducibility and easy integration into different FL frameworks.

## 2 Related Work

Numerous studies have utilized CVD data for disease detection and diagnostic support, with a particular focus on Electrocardiogram (ECG) and Echocardiogram (ECHO) data. ECG, recorded as time-series signals, captures the heart’s electrical activity over time, providing detailed insights into cardiac conditions and potential damage [1]. These studies typically formulate ECG-based tasks as classification problems aimed at disease diagnosis and heart metric analysis [20; 21; 22; 23; 24]. ECHO, which consists of ultrasound images of the heart, offers real-time visualization of heart chambers and blood flow, aiding in the diagnosis of conditions such as heart valve problems and Congestive Heart Failure (CHF). For example, [17] used ECHO data for Hypertrophic Cardiomyopathy detection through classification, while [25] employed Convolutional Neural Networks (CNNs) for standard view classification to improve clinical efficiency. Additionally, ECHO segmentation plays a crucial role in assessing heart morphological changes, and supporting diagnoses of conditions like myocardial infarction. Given that manual segmentation is time-consuming and prone to subjectivity, many studies have employed AI models for automated segmentation of ECHO images, including ventricular segmentation [26; 27; 28; 29] and atrial segmentation [30; 31]. While these studies demonstrate the potential of machine learning models on CVD data, they are largely confined to single-institution settings.

CVDs often require multi-center collaboration for effective research, and FL offers a promising solution to this challenge. Most current studies simulate FL in multi-institution collaborative training by manually partitioning data from a single

institution. For instance, [32] trained a classification model to detect cardiac arrhythmia using ECG data within a federated architecture, partitioning data from the MIT-BIH Supraventricular Arrhythmia database [33]. Similarly, [34] investigated congestive heart failure detection in a federated setting by splitting samples from the NSR-RR-interval and CHF-RR-interval databases [35] into 2 to 4 clients for simulated training. FedCluster [36] tackled the issue of unbalanced class distributions in ECG data by optimizing algorithms that cluster local parameters before performing intra- and inter-class aggregation, thus increasing the weight of minority classes. Their data were also partitioned from the MIT-BIH Arrhythmia database [35]. However, these partition-based simulations may not fully capture the true distribution characteristics of CVDs. In contrast, FLHCD [17] demonstrated federated training for hypertrophic cardiomyopathy detection using ECG and ECHO data from four medical institutions (three in the US and one in Japan), showcasing the effectiveness of FL in a naturally partitioned, multi-center setting.

To further support research in FL, numerous datasets and benchmarks have been proposed for a wide range of applications. A comprehensive comparison of FedCVD with these benchmarks is shown in Table 1. Existing studies often manually partition centralized datasets and introduce perturbations or masking to features and labels to mimic the heterogeneity found in real-world FL scenarios. For instance, FedTD [13] and FedDTI [12] simulate non-iid data partitioning by altering feature distributions and sample sizes.

To better reflect real-world conditions, several FL benchmarks utilize real-world multi-institution data directly. For instance, NIPD [15] employs data from cameras in different geographical locations as FL clients for person detection tasks, naturally exhibiting non-iid characteristics. Similarly, FEDLEGAL [16] provides a FL benchmark for NLP tasks in the legal domain, using geographically distributed case-based text data for natural data partitioning. Another example is FLMby [14], an FL benchmark for real-world distributed medical data, offering seven datasets naturally distributed by geography or institution, with corresponding tasks including segmentation and binary/multiclass classification for medical image analysis and diagnostic assistance. Some benchmarks combine natural partitioning with simulated partitioning. For instance, FedAudio [19] applies simulated partitioning for certain audio data, while introducing perturbations to mimic noisy data and labels. FedMultimodal [18] uses a Dirichlet distribution to partition multimodal data from various domains, incorporating missing modalities, labels, and erroneous labels to simulate real-world heterogeneity. Despite these advances, none of these benchmarks cover the CVD domain. Although FLHCD [17], which utilizes multi-institution data for hypertrophic cardiomyopathy detection, has a setup most similar to ours, it does not address challenges such as the long-tail distribution and incomplete label issues, which are specifically tackled by FedCVD.

### 3 The Proposed FedCVD

In this section, we present the details of the proposed general FL framework for healthcare tasks as shown in Figure 1. Our framework is built upon the lightweight open-source framework FedLab [37] for FL simulation. We present the details of datasets, metrics, and baseline models in Section 3.1. Then, we discuss the main FL challenges that FedCVD supported in Section 3.2.

#### 3.1 Datasets, Metrics and Baseline Models

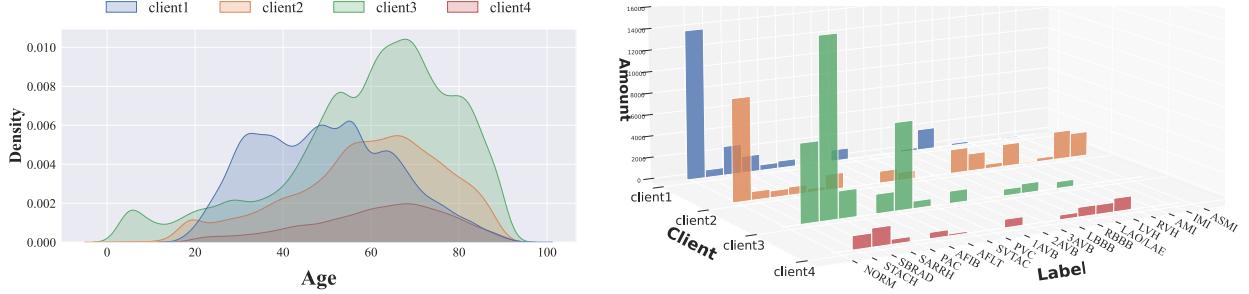
Figure 1 provides an overview of the datasets included in FedCVD. In this section, we provide a brief description of each dataset, corresponding metrics, and baseline models.

**Fed-ECG.** The 12-lead ECG signals in Fed-ECG are sourced from four distinct datasets. The first and third datasets were collected from Shandong Provincial Hospital [38] and Shaoxing People’s Hospital [39] in China, respectively. The second dataset is from the PTB-XL database, released by Physikalisch Technische Bundesanstalt (PTB) [40], and the fourth originates from the PhysioNet/Computing in Cardiology Challenge 2020 [41], which represents a large population from the Southeastern United States. These four datasets, originating from geographically diverse regions, are naturally suited for the Federated Learning (FL) setting due to their separation by location.

The corresponding task on these four datasets involves multi-label classification for each institution, a challenging problem due to the large number of labels and the long-tail distribution inherent to the data. To provide a more fine-grained evaluation, we focus on detailed label distinctions, which are of particular interest to clinicians, rather than broader label categories. To thoroughly assess performance, we adopt two metrics: *micro-F1*, which evaluates the overall performance across all labels, and *mean average precision (mAP)*, which specifically measures the impact of the long-tail distribution on model performance.

The original four datasets consist of ECG data with varying lengths and labels, each based on different standards, such as AHA [42], SCP-ECG, and SNOMED-CT, making them incompatible for use in a Federated Learning (FL) setting

directly. To standardize ECG lengths, we truncate signals longer than 5000 samples and apply edge padding to those shorter than 5000. Additionally, we retain only samples whose labels appear in at least two datasets, ensuring alignment across labels. Figures 2(a) and 2(b) illustrate the heterogeneity in both age and label distributions among institutions. For our baseline model, we adopt a residual network, following the implementation from [43]. Appendix C provides further details on the dataset and the preprocessing pipeline.



(a) Feature non-IID of the Fed-ECG dataset, demonstrated with the non-identical distribution of patient age among institutions.

(b) Label non-IID of the Fed-ECG dataset, shown as the variation in the number of each label (right axis) across different institutions (left axis).

Figure 2: Demonstration of the non-IID nature of Fed-ECG Dataset.

**Fed-ECHO.** This dataset is derived from three sources: CAMUS [44], ECHO-DYNAMIC [45], and HMC-QU [46]. CAMUS provides a database of gray-scale 2D apical four-chamber echocardiographic images, acquired at the University Hospital of St. Etienne in France, fully annotated with the left ventricular endocardium ( $LV_{Endo}$ ), epicardium ( $LV_{Epi}$ ), and left atrial wall (LA) regions. ECHO-DYNAMIC contains 2D echocardiogram videos collected at Stanford Medicine, with only the  $LV_{Endo}$  region annotated in two frames. HMC-QU, released through a collaboration between Qatar University (QU) and Hamad Medical Corporation (HMC) Hospital, includes 2D echocardiogram videos from Qatar, with annotations limited to the  $LV_{Epi}$  region in frames of a single cardiac cycle.

The common task across these three datasets is the automatic segmentation of cardiac structures in echocardiograms, a crucial step in further diagnosing cardiovascular diseases. This task is particularly challenging due to the varying quality of the original echocardiograms across datasets. To evaluate segmentation accuracy, we use both the *Dice similarity index (DICE)* and the *2D Hausdorff distance ( $d_H$ )*. The Dice index measures the overlap between the predicted segmentation and the ground truth, while the  $d_H$  quantifies the local maximum distance between the two areas.

For consistency among each institution, we only select annotated frames for experiment. The followed image preprocessing pipeline includes picture resizing to  $1 \times 112 \times 112$ , and label alignment. As a baseline model, we use a 2D U-net model following the implementation from [47]. More details about this dataset are available in Appendix D.

### 3.2 Challenging Traits of FedCVD

**Non-IID.** Non-independently and identically distributed (non-IID) is a typical characteristic in FL scenarios, encompassing non-IID features and non-IID labels, where clients' data shows heterogeneity in both feature and label spaces. Quantity imbalance, where institutions hold uneven sample sizes, can further exacerbate these non-IID issues. Among these, non-IID labels have the most pronounced impact on FL model performance. This is because the quantity and types of labels held by each institution can vary greatly, misleading the local supervised training process and causing "Client Drift" [48], which hinders global model convergence.

Fed-ECG naturally exhibits these three characteristics. In terms of feature distribution, Figure 2(a) shows significant age distribution differences among institutions' patients, with Institution 1 notably younger, reflected in the ECG features. Regarding sample size, Figure 2(b) depicts significant differences among the four institutions, with Institution 4 having the fewest samples. For label distribution in the Fed-ECG multi-label classification task, each sample may belong to multiple categories, but the quantity and proportion of different labels vary significantly among institutions. For example, the most common label for Institution 1 and Institution 2 is NORM (Normal), while for Institution 3 and Institution 4 it is STACH (Sinus tachycardia). Some institutions may even lack samples with certain labels, such as both Institution 3 and Institution 4 lacking samples labeled as PAC (Atrial premature complex(es)). These non-IID

characteristics challenge the four institutions in collaboratively training a multi-label prediction model, as institutions struggle to capture information about the labels they lack during local training, potentially leading to client drift.

**Long-tail Distribution.** In addition to the inter-institution heterogeneity caused by non-IID labels, Fed-ECG also exhibits intra-institution and inter-institution heterogeneity in the form of a long-tail distribution of labels. As shown in Figure 2(b), each institution’s internal label distribution has a clear long-tail characteristic, with a few dominant labels having many samples and numerous labels having fewer samples (long tail). These tail categories are already troublesome during independent local training, as the model may focus more on the head categories and neglect the tail ones. In FL scenarios with quantity imbalance and non-IID labels, the long-tail problem is further exacerbated. For instance, categories mainly found in the disadvantaged institutions’ tails may be in an even worse position within the overall data of all institutions. The long-tail characteristic challenges FL algorithms in ensuring the effectiveness and fairness of handling samples from various categories across institutions.

**Label Incompleteness.** Fed-ECHO presents the most challenging scenario: label-incomplete FL. In Fed-ECHO, three naturally formed institutions hold ECHO video data with annotations (image region segmentation). However, due to varying annotation capabilities, the completeness of labels among the three institutions differs, as shown in Figure 1. Institution 1 has the most complete labels (four labels) due to its ability to identify and annotate all four key regions (including the background). In contrast, Institution 2 and Institution 3 each have labels for only one key region ( $LV_{Endo}$  and  $LV_{Epi}$ , respectively). This incompleteness introduces (1) label heterogeneity, similar to the label-non-IID in Fed-ECG, where Institution 2 and Institution 3 lack some labels, and (2) mislabeling, where Institution 2 and Institution 3 label unrecognized parts as background, conflicting with Institution 1’s labels and causing misleading information. This scenario significantly challenges FL algorithms to effectively utilize the different levels of label completeness from each Institution and leverage highly heterogeneous data to benefit the global model.

## 4 Experiment

### 4.1 Experiment Details

**Baseline Algorithms.** Our experiments utilize seven typical FL algorithms across both datasets. The first four are classical global FL algorithms: *FedAvg* [9], the oft-cited FL algorithm, collaboratively trains a global model across participants. *FedProx* [49] addresses statistical heterogeneity in FL by introducing an L2 proximal term during local training, while *Scaffold* [48] mitigates client drift through control variates and server-side learning rate adjustments. *FedInit* [50] also tackles client drift by employing a personalized, relaxed initialization at the start of each local training stage. The last three are personalized FL methods: *Ditto* [51], which excels in balancing accuracy, fairness, and robustness in FL; *FedSM* [52], which combines model selection with personalized methods to avoid client drift; and *FedALA* [53], which reduces the impact of statistical heterogeneity by adaptively aggregating both the global and local models. For the Fed-ECHO dataset, we further evaluate two Federated Semi-Supervised Learning (FSSL) methods: *Fed-Consist* [10], which uses a consistency-based method for segmentation, and *FedPSL* [54], which applies separate model aggregation and meta-learning techniques for classification. In addition to the FL family, we include two other baseline algorithms: *Client*, which refers to training models using only local data without collaboration among participants, and *Central.*, which represents the ideal centralized training scenario where the server has access to all participants’ data.

**Setup.** The number of institutions involved in federated training for each task is listed in Appendix C. Our experiments mainly focus on the multi-center FL scenario (i.e., cross-silo), where all institutions participate in training at each communication round. Considering the trade-off between computation and communication, we set the local training epoch to 1 and the communication rounds to 50 throughout experiments except Fed-Consist. Since Fed-Consist requires extra rounds for training on clients with full labels before starting federated learning, we set the communication rounds to 100, where 50 rounds are for labeled clients training and another 50 rounds are normal FL training.

**Evaluation Strategies.** For a comprehensive evaluation, we build a local and global evaluation set for both datasets in FedCVD. For the local one, we divide each local data into train/test sets by 8:2. For the global one, we collect each local test set together. Our experiments test all algorithms using two evaluation strategies: 1) Global test performance (GLOBAL) is evaluated on the global test set and used to determine whether the model has learned knowledge from other clients in the FL setting. The better results of GLOBAL indicate that the model is closer to the centralized training. 2) Local test performance (LOCAL) is evaluated on each local test set. The LOCAL is more practical in real-world applications than GLOBAL because it indicates performance improvement for its task without centralizing all local data.

Table 2: The performance of different FL methods on Fed-ECG is reported using two metrics: Micro F1-Score (Mi-F1) and Mean Average Precision (mAP), both expressed as percentages (%). The best results for each configuration are highlighted in **bold**, while the second-best results are underlined.

Methods	LOCAL								GLOBAL	
	Client1		Client2		Client3		Client4		Mi-F1↑	mAP↑
	Mi-F1↑	mAP↑	Mi-F1↑	mAP↑	Mi-F1↑	mAP↑	Mi-F1↑	mAP↑		
Client1	85.8±1.9	58.1±2.6	52.7±3.4	37.8±2.2	61.5±1.2	19.8±1.2	49.8±4.2	26.7±3.0	64.3±2.1	32.3±2.0
Client2	69.9±50.0	38.9±30.0	76.8±90.0	55.7±50.0	26.3±80.0	22.7±30.0	42.2±80.0	31.6±60.0	50.4±30.0	35.9±70.0
Client3	22.7±0.2	29.8±0.7	17.0±0.4	27.2±0.3	88.1±0.2	37.7±0.4	56.9±0.4	29.4±0.6	51.5±0.2	32.7±0.2
Client4	23.7±2.0	31.7±2.7	24.7±3.3	30.5±1.5	61.6±5.5	25.3±2.1	72.3±10.2	38.5±2.8	44.7±4.3	29.3±2.5
FedAvg	69.0±10.1	58.5±1.2	50.3±5.3	54.4±0.5	77.6±0.7	37.2±0.3	66.3±0.9	39.5±0.5	67.9±3.8	50.8±0.4
FedProx	74.0±7.5	60.3±2.9	55.6±2.7	56.4±0.6	73.2±1.0	36.0±0.8	70.2±2.3	<b>43.8±1.8</b>	68.8±2.6	<b>52.3±0.9</b>
Scaffold	77.5±2.6	58.0±1.2	56.9±1.7	55.9±0.7	73.3±1.0	36.2±0.6	<u>70.7±2.9</u>	42.7±1.1	<b>70.1±0.8</b>	52.1±0.7
FedInit	73.0±6.6	58.2±0.7	54.1±5.2	55.6±1.3	73.5±0.5	36.6±0.1	67.8±2.0	41.5±1.0	68.1±3.0	51.5±0.9
Ditto	<u>82.8±4.4</u>	<b>63.1±4.2</b>	<b>74.8±1.4</b>	<b>58.3±0.6</b>	<u>86.5±1.5</u>	<b>38.1±0.6</b>	<b>73.4±6.7</b>	42.2±4.0	68.1±2.9	48.7±1.4
FedSM	77.2±7.2	58.8±1.3	59.1±4.5	56.4±1.4	69.8±0.8	35.0±0.5	67.7±3.6	42.9±2.4	68.9±2.5	51.2±0.7
FedALA	<b>84.4±4.0</b>	62.0±7.0	71.7±5.7	57.1±2.2	<b>88.2±0.1</b>	37.4±0.2	66.7±5.9	41.2±2.3	67.8±1.9	50.8±1.3
Central.	84.9±0.5	54.8±0.5	71.4±5.0	55.2±2.9	84.1±1.6	36.5±1.1	72.2±3.7	41.5±1.3	80.0±2.1	63.2±2.8

Table 3: Demonstration of FedECG’s *long-tail challenge* through the performance (F1-Score (%)) differences (measured by relative performance drop) on head and tail class groups of varying sizes. “Top K” denotes the selection of K classes with the most/fewest samples as the head/tail group. Comparisons are made among various FL algorithms, with the algorithm achieving the best result (minimum drop) highlighted in **bold** and the second-best results underlined.

Method	Top k F1 Performance											
	k=1(5%)			k=3(15%)			k=5(25%)			k=10(50%)		
	Head	Tail	Drop	Head	Tail	Drop	Head	Tail	Drop	Head	Tail	Drop
FedAvg	71.9±12.5	38.1±5.0	47.0	85.8±4.2	16.0±2.2	81.3	<b>74.1±2.5</b>	<b>25.5±1.2</b>	<b>65.6</b>	57.7±2.5	28.7±1.1	50.3
FedProx	76.8±4.5	30.8±2.6	59.9	87.6±1.4	13.1±1.2	85.0	71.2±0.9	22.7±2.6	68.1	57.5±1.8	31.1±2.2	46.0
Scaffold	77.4±3.2	33.8±4.9	56.4	87.7±1.2	14.6±2.2	83.4	71.3±0.7	24.1±1.4	<b>66.3</b>	<b>58.4±1.7</b>	<b>32.1±2.3</b>	<b>45.0</b>
FedInit	77.2±2.8	29.7±2.7	61.5	87.8±1.0	12.9±1.2	87.9	71.3±0.6	23.2±0.5	69.8	56.4±2.5	29.0±0.8	45.8
Ditto	73.5±7.4	23.8±5.8	67.6	86.4±2.5	10.2±1.7	88.2	70.6±2.1	21.4±1.1	69.7	54.4±2.6	27.6±1.7	49.3
FedSM	75.5±8.2	25.8±3.6	65.8	87.2±2.9	10.5±2.2	85.3	69.5±1.8	21.0±0.9	67.5	57.0±1.9	30.9±2.3	48.5
FedALA	<b>72.4±5.4</b>	<b>38.9±5.8</b>	<b>46.2</b>	<b>86.0±1.8</b>	<b>16.2±2.1</b>	<b>81.1</b>	73.7±1.5	25.2±1.2	65.7	57.9±1.7	28.8±1.0	50.3
Central.	88.6±2.3	35.6±5.9	59.8	92.4±1.6	19.5±7.0	78.9	84.1±1.9	29.9±5.5	64.5	71.3±3.6	44.5±4.4	37.6

## 4.2 Benchmark on Fed-ECG

We present the Fed-ECG dataset, which poses significant challenges for FL scenarios, namely *non-IID data* and *long-tailed distribution*. We first compared the overall performance of mainstream FL algorithms on Fed-ECG, with the evaluated local and global performance shown in Table 2. The results indicate that FL has advantages over local training. Additionally, FL algorithms designed for heterogeneous scenarios (FedProx, Scaffold, FedInit) outperform the FedAvg algorithm. The personalized algorithms Ditto, FedSM, and FedALA exhibit excellent performance in local tests for certain clients. However, these FL algorithms still lag behind centralized training.

To better illustrate the impact of these two challenges on FL performance, we introduced additional experimental settings and evaluation metrics. For the *non-IID challenge*, we compared the performance differences between natural partitioning and two simulated partitioning methods (random and non-IID), with the simulated non-IID partitioning method described in Appendix C. Figure 3 compares the performance (percentage relative to centralized training) between FL algorithms trained under the three partitioning settings. The results reveal that Fed-ECG’s natural partitioning poses significantly greater challenges compared to the two simulated partitioning methods.

For the *long-tail challenge*, we used the mAP metric in Table 2 to evaluate the overall performance of algorithms across different classes. In general, FL algorithms designed for heterogeneous scenarios demonstrate an advantage in addressing long-tail issues, with personalized algorithms like Ditto and FedALA showing better results in local tests. However, in global tests, the FedProx algorithm more effectively handles long-tail problems. Comparisons with centralized training reveal that FL scenarios tend to amplify the impact of long-tail distributions. To further illustrate this challenge, we introduced the F1-STD metric, which measures the standard deviation of F1 scores across classes. This metric reflects the algorithm’s ability to manage long-tail problems; the larger the STD, the poorer the algorithm’s performance in this regard. The GLOBAL F1-STD results of different FL algorithms are visually presented in Figure 4, showing a pattern consistent with Table 2 and underscoring the challenges posed by long-tail distributions.

To more granularly evaluate the performance of algorithms in long-tail scenarios, we introduced a new metric, *Top-K*. Top-K refers to selecting the K classes with the most samples and the K classes with the fewest samples, calculating the average F1 score for each group, and then computing the relative performance drop between them. A larger performance

drop indicates a more severe long-tail problem. Table 3 presents the Top-K metrics for various K values, highlighting the significant long-tail characteristics of Fed-ECG. The results show that mainstream FL algorithms struggle to address long-tail issues effectively, performing worse compared to centralized training.

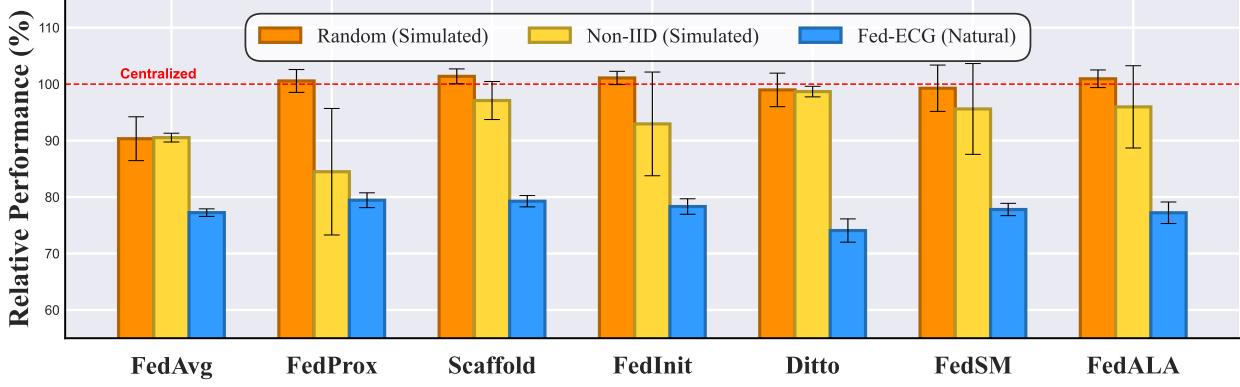


Figure 3: Demonstration of Fed-ECG’s *non-IID challenge*: Comparisons of performance (relative Mean Average Score %) between artificial partitions (simulated random and non-IID partitions) and Fed-ECG’s natural partition across different FL algorithms.

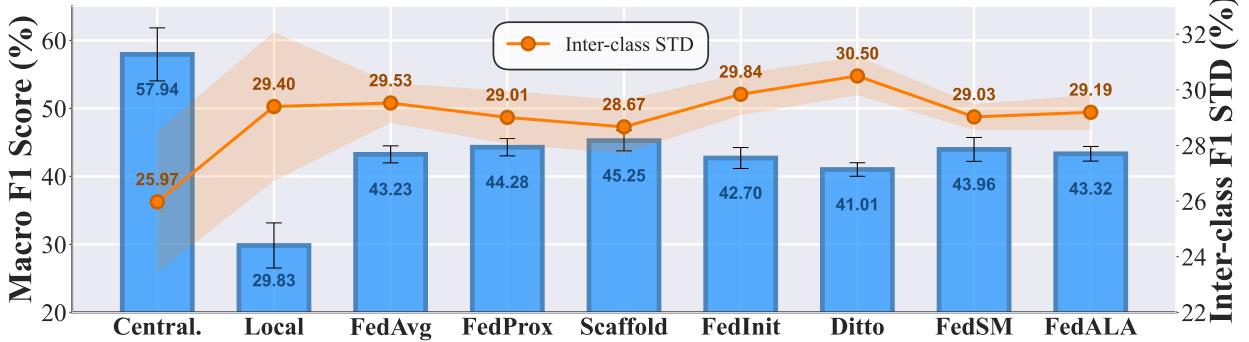


Figure 4: Demonstration of Fed-ECG’s *long-tail challenge*: Average Macro F1-Score (%) and Standard Deviation across classes for various FL Algorithms.

### 4.3 Benchmark on Fed-ECHO

The proposed Fed-ECHO dataset presents one of the most challenging FL settings: *label incompleteness*, which can be viewed as an enhanced version of label-non-IID. Specifically, all annotated video frames from Institution 1 are completely segmented into four regions: BG, LV<sub>Endo</sub>, LV<sub>Epi</sub>, and LA. In contrast, Institution 2 and Institution 3 can only recognize the LV<sub>Epi</sub> and LV<sub>Endo</sub> regions in their annotated video frames, respectively, with the remaining regions simply labeled as "BG." For convenience, we refer to the BG labels from Institution 2 and Institution 3 as "Maybe-BG," indicating these segmentations may be unreliable. This discrepancy introduces potential conflicts between the "Maybe-BG" labels of Institution 2 and Institution 3 and the corresponding "reliable" labels of Institution 1, resulting in misleading labels that affect model convergence.

To mitigate the impact of misleading labels, we propose a straightforward baseline strategy, *supervised-only*. During supervised model training, we input the data from all three Institutions into the model without additional processing, allowing the model to benefit from the rich data features. However, when calculating the loss, we mask out the "Maybe-BG" regions in the video frames from Institution 2 and Institution 3. This means that for samples from Institution 2 and Institution 3, we only compute the training loss on the "reliable foreground". This strategy ensures the model learns segmentation capabilities from completely reliable labels. Additionally, during model segmentation performance evaluation, we also exclude the "Maybe-BG" regions from the test samples, preventing them from influencing the model’s performance metrics.

Table 4: The performance of different FL methods on Fed-ECHO, with DICE (%) and  $d_H$  representing DICE index and Hausdorff distance respectively. The best results for each configuration are highlighted in **bold**, while the second-best results are underlined.

Mthods	LOCAL						GLOBAL	
	Client1		Client2		Client3		Dice↑	$d_H \downarrow$
	Dice↑	$d_H \downarrow$	Dice↑	$d_H \downarrow$	Dice↑	$d_H \downarrow$		
Client1	88.2±0.8	5.196±0.360	46.5±3.9	24.246±0.442	63.4±4.2	22.000±12.914	66.1±2.8	17.147±4.187
Client2	24.4±6.0	71.917±2.832	88.9±5.5	5.577±1.413	-	-	37.8±3.8	59.165±1.411
Client3	15.8±1.0	76.368±0.988	-	-	94.1±0.7	7.110±2.900	36.6±0.3	61.159±0.931
FedAvg	26.2±3.7	48.343±8.719	56.4±8.8	33.127±10.721	67.9±3.5	34.004±5.287	50.2±5.3	38.491±8.058
FedProx	74.8±18.7	13.928±11.742	<b>82.3±3.5</b>	13.402±2.975	66.7±12.4	16.181±16.329	74.6±11.4	14.504±9.134
Scaffold	81.5±2.1	9.981±2.482	81.0±2.1	<u>12.543±2.157</u>	<b>74.6±2.1</b>	<u>7.551±0.885</u>	<u>79.0±0.7</u>	<u>10.025±1.467</u>
FedInit	83.5±0.9	<u>7.799±0.665</u>	<u>81.6±2.2</u>	<b>12.240±1.091</b>	73.4±3.0	<b>7.542±0.918</b>	<b>79.5±0.5</b>	<b>9.193±0.558</b>
Ditto	<b>88.2±0.4</b>	<b>4.796±0.085</b>	56.9±3.3	28.381±4.043	56.3±2.2	27.321±15.627	78.1±1.8	10.658±2.372
FedSM	80.2±6.0	11.339±5.868	81.1±1.5	12.580±1.288	72.7±2.0	10.913±4.128	78.0±2.2	11.611±2.308
FedALA	80.5±1.6	8.700±1.245	51.3±2.4	36.472±2.686	47.1±0.9	52.128±4.356	52.3±2.0	36.811±2.630
Fed-Consist	85.9±0.2	11.904±0.442	75.2±0.9	27.480±1.440	66.3±0.2	34.037±1.777	75.8±0.3	24.474±1.155
FedPSL	53.5±9.3	37.277±9.166	77.0±2.9	12.873±1.589	67.8±14.1	29.166±15.660	66.1±7.5	26.439±7.831
Central.(sup)	89.9±0.4	4.643±0.097	48.5±22.2	43.684±19.659	65.0±14.6	30.557±14.831	67.8±12.1	26.295±11.379
Central.(ssup)	90.3±0.2	3.872±0.067	91.7±0.5	4.370±0.181	91.1±1.7	3.005±0.732	91.0±0.6	3.749±0.242

Table 4 compares the performance of mainstream FL algorithms with centralized/isolated learning on Fed-ECHO, evaluated using Dice and Hausdorff distance ( $d_H$ ). Except for the semi-supervised learning algorithms Centralized (semi-sup) and Fed-Consist, all algorithms use the previously mentioned supervised-only strategy. The results underscore the viability of FL in the Fed-ECHO setting, as most FL algorithms exhibit superior global performance compared to models trained independently by individual institutions. However, due to high degree of data heterogeneity, none of the evaluated FL algorithms outperform locally trained models on each client’s test dataset, indicating a need for more personalized and heterogeneity-resistant FL strategies.

On the global test set, FL algorithms specifically designed to address heterogeneity, such as FedInit and Scaffold, consistently demonstrate significant advantages over simpler algorithms like FedAvg. Notably, these algorithms also outperform the Centralized (sup) model, which we attribute to FL effectively mitigating the impact of intra-batch heterogeneity(e.g., within the same batch, there are four-label data from Institution 1 and single-label data from Institution 2 or Institution 3).

Additionally, to leverage the substantial amount of partially labeled data from Institution 2 and Institution 3 and potentially mitigate label heterogeneity, we introduced semi-supervised learning algorithms for comparison. These include the centralized semi-supervised model, Centralized (ssup), and federated semi-supervised algorithms such as Fed-Consist and FedPSL. The Centralized (ssup) model significantly outperformed its fully supervised counterpart, underscoring the value of utilizing unlabeled video frames. Similarly, Fed-Consist outperformed FedAvg and FedProx, although it still exhibited a noticeable performance gap compared to the centralized semi-supervised algorithm and lagged behind fully supervised FL algorithms like Scaffold and FedInit. While FedPSL performed well on certain participating client, it showed greater instability overall, largely due to its sensitivity to client-side feature heterogeneity.

Therefore, the highly heterogeneous Fed-ECHO scenario poses significant challenges for FL algorithms, requiring them to adapt to heterogeneous data and effectively leverage unlabeled data across different data domains.

## 5 Conclusion

This paper has introduced FedCVD, the first real-world multi-center FL benchmark for CVD data, which consists of two datasets and their respective tasks: Fed-ECG and Fed-ECHO. It presents three major challenges due to the heterogeneous distribution of real-world data: non-IID, long-tailed labels, and label incompleteness. We conducted extensive comparative and validation experiments, testing mainstream FL algorithms and centralized training on these tasks. Experimental results show that the natural non-IID characteristics in FedCVD are more challenging than the manually partitioned setups in most previous federated benchmarks, and mainstream algorithms perform poorly in the long-tail tests of FedCVD. For the most difficult task, i.e., the label-incomplete Fed-ECHO, mainstream FL algorithms barely maintain utility, but are still better than non-cooperative algorithms that only utilize unlabeled data on each client. Federated semi-supervised learning algorithms that leverage unlabeled data achieve some performance improvement. Beyond, as a flexible and extensible framework, FedCVD is meant to be a step towards the development of FL on CVD domain.

**Limitations and Future Work.** FedCVD presents a realistic and challenging scenario that tests FL algorithms' ability to mitigate data heterogeneity, handle long-tailed classes, and utilize unlabeled data. However, FedCVD currently offers a limited variety of data types and only two tasks. Additionally, the FL algorithms compared in experiments, particularly semi-supervised ones, are limited. In future work, we will expand the data range of FedCVD, aiming for it to inspire future FL research in real-world medical contexts, especially with CVD data.

## References

- [1] Sricharan Donkada, Seyedamin Pouriyeh, Reza M. Parizi, Meng Han, Nasrin Dehbozorgi, Nazmus Sakib, and Quan Z. Sheng. Uncovering promises and challenges of federated learning to detect cardiovascular diseases: A scoping literature review. *CoRR*, abs/2308.13714, 2023.
- [2] Lerina Aversano, Mario Luca Bernardi, Marta Cimitile, Martina Iammarino, Debora Montano, and Chiara Verdone. Using machine learning for early prediction of heart disease. In Plamen Angelov, George A. Papadopoulos, Giovanna Castellano, José A. Iglesias, Gabriella Casalino, Edwin Lughofer, and Daniel Leite, editors, *IEEE International Conference on Evolving and Adaptive Intelligent System, EAIS 2022, Larnaca, Cyprus, May 25-26, 2022*, pages 1–8. IEEE, 2022.
- [3] Li Yan, Hai-Tao Zhang, Jorge M. Gonçalves, Yang Xiao, Maolin Wang, Yuqi Guo, Chuan Sun, Xiuchuan Tang, Liang Jing, Mingyang Zhang, Xiang Huang, Ying Xiao, Haosen Cao, Yanyan Chen, Tongxin Ren, Fang Wang, Yaru Xiao, Sufang Huang, Xi Tan, Niannian Huang, Bo Jiao, Cheng Cheng, Yong Zhang, Ailin Luo, Laurent Mombaerts, Junyang Jin, Zhiguo Cao, Shusheng Li, Hui Xu, and Ye Yuan. An interpretable mortality prediction model for COVID-19 patients. *Nat. Mach. Intell.*, 2(5):283–288, 2020.
- [4] Lingxiao Chen. Overview of clinical prediction models. *Annals of translational medicine*, 8(4), 2020.
- [5] Roohallah Alizadehsani, Moloud Abdar, Mohamad Roshanzamir, Abbas Khosravi, Parham M. Kebria, Fahime Khozeimeh, Saeid Nahavandi, Nizal Sarratzadeegan, and U. Rajendra Acharya. Machine learning-based coronary artery disease diagnosis: A comprehensive review. *Comput. Biol. Medicine*, 111, 2019.
- [6] Subhi J Al’Aref, Khalil Anchouche, Gurpreet Singh, Piotr J Slomka, Kranthi K Kolli, Amit Kumar, Mohit Pandey, Gabriel Maliakal, Alexander R Van Rosendael, Ashley N Beecy, et al. Clinical applications of machine learning in cardiovascular disease and its relevance to cardiac imaging. *European heart journal*, 40(24):1975–1986, 2019.
- [7] Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA relevance), April 2016. Legislative Body: EP, CONSIL.
- [8] Qiang Yang, Yang Liu, Yong Cheng, Yan Kang, Tianjian Chen, and Han Yu. *Federated Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2019.
- [9] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In Aarti Singh and Xiaojin (Jerry) Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 2017.
- [10] Dong Yang, Ziyue Xu, Wenqi Li, Andriy Myronenko, Holger R Roth, Stephanie Harmon, Sheng Xu, Baris Turkbey, Evrim Turkbey, Xiaosong Wang, et al. Federated semi-supervised learning for covid region segmentation in chest ct using multi-national data from china, italy, japan. *Medical image analysis*, 70:101992, 2021.
- [11] Karthik V Sarma, Stephanie Harmon, Thomas Sanford, Holger R Roth, Ziyue Xu, Jesse Tetreault, Daguang Xu, Mona G Flores, Alex G Raman, Rushikesh Kulkarni, et al. Federated learning improves site performance in multicenter deep learning without data sharing. *Journal of the American Medical Informatics Association*, 28(6):1259–1264, 2021.
- [12] Gianluca Mittone, Filip Svoboda, Marco Aldinucci, Nicholas D. Lane, and Pietro Lió. A federated learning benchmark for drug-target interaction. In Ying Ding, Jie Tang, Juan F. Sequeda, Lora Aroyo, Carlos Castillo, and Geert-Jan Houben, editors, *Companion Proceedings of the ACM Web Conference 2023, WWW 2023, Austin, TX, USA, 30 April 2023 - 4 May 2023*, pages 1177–1181. ACM, 2023.
- [13] William Lindskog and Christian Prehofer. A federated learning benchmark on tabular data: Comparing tree-based models and neural networks. In *Eighth International Conference on Fog and Mobile Edge Computing, FMEC 2023, Tartu, Estonia, September 18-20, 2023*, pages 239–246. IEEE, 2023.
- [14] Jean Ogier du Terrail, Samy-Safwan Ayed, Edwige Cyffers, Felix Grimberg, Chaoyang He, Regis Loeb, Paul Mangold, Tanguy Marchand, Othmane Marfoq, Erum Mushtaq, Boris Muzellec, Constantin Philippenko, Santiago Silva, Maria Telenczuk, Shadi Albarqouni, Salman Avestimehr, Aurélien Bellet, Aymeric Dieuleveut, Martin Jaggi, Sai Praneeth Karimireddy, Marco Lorenzi, Giovanni Neglia, Marc Tommasi, and Mathieu Andreux. Flambý: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.

- [15] Kangning Yin, Zhen Ding, Zhihua Dong, Dongsheng Chen, Jie Fu, Xinhui Ji, Guangqiang Yin, and Zhiguo Wang. NIPD: A federated learning person detection benchmark based on real-world non-iid data. *CoRR*, abs/2306.15932, 2023.
- [16] Zhuo Zhang, Xiangjing Hu, Jingyuan Zhang, Yating Zhang, Hui Wang, Lizhen Qu, and Zenglin Xu. FEDLEGAL: the first real-world federated learning benchmark for legal NLP. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 3492–3507. Association for Computational Linguistics, 2023.
- [17] Shinichi Goto, Divyarajsinhji Solanki, Jenine E John, Ryuichiro Yagi, Max Homilius, Genki Ichihara, Yoshinori Katsumata, Hanna K Gaggin, Yuji Itabashi, Calum A MacRae, et al. Multinational federated learning approach to train ecg and echocardiogram models for hypertrophic cardiomyopathy detection. *Circulation*, 146(10):755–769, 2022.
- [18] Tiantian Feng, Digbalay Bose, Tuo Zhang, Rajat Hebbar, Anil Ramakrishna, Rahul Gupta, Mi Zhang, Salman Avestimehr, and Shrikanth Narayanan. Fedmultimodal: A benchmark for multimodal federated learning. In Ambuj K. Singh, Yizhou Sun, Leman Akoglu, Dimitrios Gunopulos, Xifeng Yan, Ravi Kumar, Fatma Ozcan, and Jieping Ye, editors, *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2023, Long Beach, CA, USA, August 6-10, 2023*, pages 4035–4045. ACM, 2023.
- [19] Tuo Zhang, Tiantian Feng, Samiul Alam, Sunwoo Lee, Mi Zhang, Shrikanth S. Narayanan, and Salman Avestimehr. Fedaudio: A federated learning benchmark for audio tasks. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE, 2023.
- [20] Dusit Thanapatay, Chaiwat Suwansaroj, and Chusak Thanawattano. Ecg beat classification method for ecg printout with principle components analysis and support vector machines. In *2010 International Conference on Electronics and Information Engineering*, volume 1, pages V1–72. IEEE, 2010.
- [21] RJ Muirhead and RD Puff. A bayesian classification of heart rate variability data. *Physica A: Statistical Mechanics and its Applications*, 336(3-4):503–513, 2004.
- [22] I Christov, I Jekova, and G Bortolan. Premature ventricular contraction classification by the kth nearest-neighbours rule. *Physiological measurement*, 26(1):123, 2005.
- [23] Omar Behadada and Mohammed Amine Chikh. An interpretable classifier for detection of cardiac arrhythmias by using the fuzzy decision tree. *Artif. Intell. Res.*, 2(3):45–58, 2013.
- [24] Jing Zhang, Deng Liang, Aiping Liu, Min Gao, Xiang Chen, Xu Zhang, and Xun Chen. Mlbf-net: A multi-lead-branch fusion network for multi-class arrhythmia classification using 12-lead ecg. *IEEE journal of translational engineering in health and medicine*, 9:1–11, 2021.
- [25] Andreas Ostvik, Erik Smistad, Svein Arne Aase, Bjørn Olav Haugen, and Lasse Lovstakken. Real-time standard view classification in transthoracic echocardiography using convolutional neural networks. *Ultrasound in medicine & biology*, 45(2):374–384, 2019.
- [26] Yaonan Zhang, Yuan Gao, Jinling Jiao, Xian Li, Sai Li, and Jun Yang. Robust boundary detection and tracking of left ventricles on ultrasound images using active shape model and ant colony optimization. *Bio-Medical Materials and Engineering*, 24(6):2893–2899, 2014.
- [27] Auzuir Ripardo De Alexandria, Paulo César Cortez, Jessyca Almeida Bessa, John Hebert da Silva Félix, José Sebastião De Abreu, and Victor Hugo C De Albuquerque. psnakes: A new radial active contour model and its application in the segmentation of the left ventricle from echocardiographic images. *Computer methods and programs in biomedicine*, 116(3):260–273, 2014.
- [28] Xulei Qin, Zhibin Cong, and Baowei Fei. Automatic segmentation of right ventricular ultrasound images using sparse matrix transform and a level set. *Physics in Medicine & Biology*, 58(21):7609, 2013.
- [29] Serkan Kiranyaz, Aysen Degerli, Tahir Hamid, Rashid Mazhar, Rayyan El Fadil Ahmed, Rayaan Abouhasera, Morteza Zabihi, Junaid Malik, Ridha Hamila, and Moncef Gabbouj. Left ventricular wall motion estimation by active polynomials for acute myocardial infarction detection. *IEEE Access*, 8:210301–210317, 2020.
- [30] Alexander Haak, Gonzalo Vegas-Sánchez-Ferrero, Harriet W Mulder, Ben Ren, Hortense A Kirişli, Coert Metz, Gerard van Burken, Marijn van Stralen, Josien PW Pluim, FW van der Steen, et al. Segmentation of multiple heart cavities in 3-d transesophageal ultrasound images. *IEEE transactions on ultrasonics, ferroelectrics, and frequency control*, 62(6):1179–1189, 2015.
- [31] Alexander Haak, Ben Ren, Harriet W Mulder, Gonzalo Vegas-Sánchez-Ferrero, Gerard van Burken, Antonius FW van der Steen, Marijn van Stralen, Josien PW Pluim, Theo van Walsum, and Johannes G Bosch. Improved

- segmentation of multiple cavities of the heart in wide-view 3-d transesophageal echocardiograms. *Ultrasound in medicine & biology*, 41(7):1991–2000, 2015.
- [32] Sadman Sakib, Mostafa M. Fouda, Zubair Md. Fadlullah, Khalid Abualsaoud, Elias Yaacoub, and Mohsen Guizani. Asynchronous federated learning-based ECG analysis for arrhythmia detection. In *IEEE International Mediterranean Conference on Communications and Networking, MeditCom 2021, Athens, Greece, September 7-10, 2021*, pages 277–282. IEEE, 2021.
  - [33] Scott David Greenwald, Ramesh S Patil, and Roger G Mark. *Improved detection and classification of arrhythmias in noise-corrupted electrocardiograms using contextual information*. IEEE, 1990.
  - [34] Liang Zou, Zexin Huang, Xinhui Yu, Jiannan Zheng, Aiping Liu, and Meng Lei. Automatic detection of congestive heart failure based on multiscale residual unet++: From centralized learning to federated learning. *IEEE Trans. Instrum. Meas.*, 72:1–13, 2023.
  - [35] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.
  - [36] Daoqin Lin, Yuchun Guo, Huan Sun, and Yishuai Chen. Fedcluster: A federated learning framework for cross-device private ECG classification. In *IEEE INFOCOM 2022 - IEEE Conference on Computer Communications Workshops, INFOCOM 2022 - Workshops, New York, NY, USA, May 2-5, 2022*, pages 1–6. IEEE, 2022.
  - [37] Dun Zeng, Siqi Liang, Xiangjing Hu, Hui Wang, and Zenglin Xu. Fedlab: A flexible federated learning framework. *Journal of Machine Learning Research*, 24(100):1–7, 2023.
  - [38] Hui Liu, Dan Chen, Da Chen, Xiyu Zhang, Huijie Li, Lipan Bian, Minglei Shu, and Yinglong Wang. A large-scale multi-label 12-lead electrocardiogram database with standardized diagnostic statements. *Scientific data*, 9(1):272, 2022.
  - [39] J Zheng, H Guo, and H Chu. A large scale 12-lead electrocardiogram database for arrhythmia study (version 1.0. 0). *PhysioNet 2022 Available online [http://physionet.org/content/ecg\\_arrhythmia10](http://physionet.org/content/ecg_arrhythmia10) accessed on*, 23, 2022.
  - [40] Patrick Wagner, Nils Strodthoff, Ralf-Dieter Bousseljot, Dieter Kreiseler, Fatima I Lunze, Wojciech Samek, and Tobias Schaeffter. PtB-XL, a large publicly available electrocardiography dataset. *Scientific data*, 7(1):1–15, 2020.
  - [41] Erick A Perez Alday, Annie Gu, Amit J Shah, Chad Robichaux, An-Kwok Ian Wong, Chengyu Liu, Feifei Liu, Ali Bahrami Rad, Andoni Elola, Salman Seyed, et al. Classification of 12-lead ecgs: the physionet/computing in cardiology challenge 2020. *Physiological measurement*, 41(12):124003, 2020.
  - [42] Paul Kligfield, Leonard S Gettes, James J Bailey, Rory Childers, Barbara J Deal, E William Hancock, Gerard Van Herpen, Jan A Kors, Peter Macfarlane, David M Mirvis, et al. Recommendations for the standardization and interpretation of the electrocardiogram: part i: the electrocardiogram and its technology: a scientific statement from the american heart association electrocardiography and arrhythmias committee, council on clinical cardiology; the american college of cardiology foundation; and the heart rhythm society endorsed by the international society for computerized electrocardiology. *Circulation*, 115(10):1306–1324, 2007.
  - [43] Nils Strodthoff, Patrick Wagner, Tobias Schaeffter, and Wojciech Samek. Deep learning for ecg analysis: Benchmarks and insights from ptb-xl. *IEEE journal of biomedical and health informatics*, 25(5):1519–1528, 2020.
  - [44] Sarah Leclerc, Erik Smistad, Joao Pedrosa, Andreas Østvik, Frederic Cervenansky, Florian Espinosa, Torvald Espeland, Erik Andreas Rye Berg, Pierre-Marc Jodoin, Thomas Grenier, et al. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE transactions on medical imaging*, 38(9):2198–2210, 2019.
  - [45] David Ouyang, Bryan He, Amirata Ghorbani, Neal Yuan, Joseph Ebinger, Curtis P Langlotz, Paul A Heidenreich, Robert A Harrington, David H Liang, Euan A Ashley, et al. Video-based ai for beat-to-beat assessment of cardiac function. *Nature*, 580(7802):252–256, 2020.
  - [46] Serkan Kiranyaz, Aysen Degerli, Tahir Hamid, Rashid Mazhar, Rayyan El Fadil Ahmed, Rayaan Abouhasera, Morteza Zabihi, Junaid Malik, Ridha Hamila, and Moncef Gabbouj. Left ventricular wall motion estimation by active polynomials for acute myocardial infarction detection. *IEEE Access*, 8:210301–210317, 2020.
  - [47] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III* 18, pages 234–241. Springer, 2015.
  - [48] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning*, pages 5132–5143. PMLR, 2020.

- [49] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020.
- [50] Yan Sun, Li Shen, and Dacheng Tao. Understanding how consistency works in federated learning via stage-wise relaxed initialization. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [51] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *International conference on machine learning*, pages 6357–6368. PMLR, 2021.
- [52] An Xu, Wenqi Li, Pengfei Guo, Dong Yang, Holger Roth, Ali Hatamizadeh, Can Zhao, Daguang Xu, Heng Huang, and Ziyue Xu. Closing the generalization gap of cross-silo federated medical image segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 20834–20843. IEEE, 2022.
- [53] Jianqing Zhang, Yang Hua, Hao Wang, Tao Song, Zhengui Xue, Ruhui Ma, and Haibing Guan. Fedala: Adaptive local aggregation for personalized federated learning. In Brian Williams, Yiling Chen, and Jennifer Neville, editors, *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 11237–11244. AAAI Press, 2023.
- [54] Nanqing Dong, Michael Kampffmeyer, Irina Voiculescu, and Eric P. Xing. Federated partially supervised learning with limited decentralized medical images. *IEEE Trans. Medical Imaging*, 42(7):1944–1954, 2023.

## A Broader Impact

Considering that this research exclusively involves the repurposing of existing open-source databases, the associated risks are limited. However, it is important to acknowledge that all datasets utilized in this study may be influenced by biases inherent in the original data collection processes, such as those related to gender, age, or race. Unfortunately, identifying the sources of potential biases is challenging because the data have been appropriately pseudonymized. Moreover, records such as electrocardiograms and echocardiograms cannot be easily linked to specific demographic attributes such as age, ethnicity, or gender by non-medical experts. Nonetheless, our work discloses certain metadata of the datasets, including geographical origin, gender distribution, and age distribution. This exposure may aid in identifying underlying geographical biases, which are anticipated in real-world federated learning scenarios.

While prioritizing simplicity and utility, the current benchmark does not include privacy metrics. Nevertheless, privacy remains critically important in the cardiovascular disease domain, and we strongly encourage the research community to address these considerations. Thanks to the modularity of FedCVD, we can add privacy components easily. Therefore, we anticipate that FedCVD will address privacy concerns related to federated learning within the cardiovascular disease domain in the future.

## B Datasets repository and Maintenance plane

### B.1 Dataset repository.

The code is now available at <https://github.com/SMILELab-FL/FedCVD>. Considering licenses, users need to download the data manually through the original dataset link.

### B.2 Maintenance plan

We shall adhere to a maintenance plan to uphold the integrity of the codebase and ensure the conformity of supplied datasets to requisite standards. In particular, this maintenance plan encompasses:

- Fixing bugs affecting the correctness of our code, whether identified by the community or ourselves;
- Introducing additional variants of federated learning techniques, including alternative methods within the scope of cross-silo federated learning and federated semi-supervised learning methodologies;
- Adding new functional modules, such as privacy protection components.
- Regarding datasets, reviewing potential updates of the datasets referenced in the FedCVD, including but not limited to introducing new tasks or modalities;

## C Fed-ECG

### C.1 Description

Fed-ECG consists of four datasets: SPH, PTB-XL, SXPH, and G12EC. The order of leads of each dataset is I, II, III, aVR, aVL, aVF, V1, V2, V3, V4, V5, V6. The overview of Fed-ECG is shown in Table 5. Table 6 shows demographics information for four datasets in Fed-ECG.

**SPH.** The original Shandong Provincial Hospital (SPH) database contains 25,770 12-lead ECG records from 24,666 patients, which were acquired from Shandong Provincial Hospital between 2019/08 and 2020/08. The record length is between 10 and 60 seconds. The sampling frequency is 500 Hz. All ECG records are in full compliance with the AHA standard, which aims for the standardization and interpretation of the electrocardiogram and consists of 44 primary statements and 15 modifiers as per the standard. 46.04% records in this dataset contain ECG abnormalities. Moreover, 14.45% records have multiple diagnostic statements.

**PTB-XL.** The original PTB-XL database contains 21,837 12-lead ECG records from 18,885 patients of 10 seconds length at the Physikalisch Technische Bundesanstalt (PTB) between October 1989 and June 1996. The original records are resampled to both 100 Hz and 500 Hz. For consistency, we only use the records whose frequency is 500 Hz. Each data is annotated by up to two cardiologists with the SCP-ECG standard.

Table 5: Overview of the datasets, tasks, metrics and baseline models in FedCVD.

Dataset	Fed-ECG				Fed-ECHO		
Task Type	Multi-label Classification				2D Segmentation		
Input	12-lead ECG Signal				Echocardiogram		
Prediction (y)	Diagnostic Statement				Cardiac Structure Mask		
Data source	SPH	PTB-XL	SXPH	G12EC	CAMUS	ECHONET-DYNAMIC	HMC-QU
Original Patient Size	24,666	18,885	45,152	UNKNOWN	500	10,030	109
Original Sample Size	25,770	21,837	45,152	10,344	1000	20,060	2,349
Preprocessing	Label Alignment				Resizing and Label Alignment		
Patient Size	21,530	16,699	36,272	UNKNOWN	500	10,024	109
Sample Size	22,425	19,019	36,272	6,205	1000	20,048	2,349
Model	ResNet				U-net		
Metrics	Micro F1 / mAP				DICE / Hausdorff distance		
Input Dimension	$12 \times 5000$				$112 \times 112$		

**SXPH.** This database contains 12-lead ECGs of 45,152 patients with a 500 Hz sampling rate under the auspices of Chapman University, Shaoxing People’s Hospital (Shaoxing Hospital Zhejiang University School of Medicine), and Ningbo First Hospital. The record length is 10 seconds. All records are labeled by professional experts with the SNOMED-CT standard.

**G12EC.** This Georgia 12-lead ECG Challenge (G12EC) database is provided by the PhysioNet/Computing in Cardiology Challenge 2020. Only 10,344 training data from this database are open to the public. The record length is not longer than 10 seconds with a sample frequency of 500 Hz. All records are labeled with the SNOMED-CT standard as well.

Table 6: Demographics information for Fed-ECG.

Client	Sex	Dataset size	Age	Age Range
Client1	Female	9,502	$48.73 \pm 15.67$	18 - 92
	Male	12,923	$50.35 \pm 15.49$	18 - 95
Client2	Female	8,930	$59.80 \pm 18.42$	3 - 89
	Male	10,089	$58.40 \pm 15.66$	2 - 89
Client3	Female	14,830	$58.36 \pm 20.11$	4 - 89
	Male	21,442	$60.28 \pm 19.10$	4 - 89
Client4	Female	2,668	$61.37 \pm 16.51$	20 - 89
	Male	3,537	$61.35 \pm 15.04$	14 - 89

## C.2 License and Ethics

All four databases are open-access. The SPH database is open access at Figshare, while the rest databases are open access at PhysioNet under a Creative Commons Attribution 4.0 International Public License.

The PTB-XL database was supported by the Bundesministerium für Bildung und Forschung (BMBF) through the Berlin Big Data Center under Grant 01IS14013A and the Berlin Center for Machine Learning under Grant 01IS18037I and by the EMPIR project 18HLT07 MedalCare. The EMPIR initiative is cofunded by the European Union’s Horizon 2020 research and innovation program and the EMPIR Participating States.

The institutional review board of Shaoxing People’s Hospital and Ningbo First Hospital of Zhejiang University approved the study of the SXPH database, granted the waiver application to obtain informed consent, and allowed the data to be shared publicly after de-identification. The requirement for patient consent was waived.

## C.3 Download and preprocessing

### C.3.1 Download

The four datasets can be downloaded using the URLs below:

1. **SPH:** [https://springernature.figshare.com/collections/A\\_large-scale\\_multi-label\\_12-lead\\_ecg\\_database\\_with\\_standardized\\_diagnostic\\_statements/5779802/1](https://springernature.figshare.com/collections/A_large-scale_multi-label_12-lead_ecg_database_with_standardized_diagnostic_statements/5779802/1)

Table 7: Label relationship between original label and ours.

ours	SPH	Original Label		
		PTB-XL	SXPH	G12EC
NORM (Normal)	Normal	Normal	-	-
STACH (Sinus tachycardia)	Sinus tachycardia	Sinus tachycardia	Sinus tachycardia	427084000
SBRAD (Sinus bradycardia)	Sinus bradycardia	Sinus bradycardia	Sinus bradycardia	426177001
SARRH (Sinus arrhythmia)	Sinus arrhythmia	Sinus arrhythmia	-	427393009
PAC (Atrial premature complex(es))	Atrial premature complex(es)	Atrial premature complex	-	-
AFIB (Atrial fibrillation)	Atrial fibrillation	Atrial fibrillation	Atrial fibrillation	164889003
AFLT (Atrial flutter)	Atrial flutter	Atrial flutter	Atrial flutter	164890007
SVTAC (Supraventricular tachycardia)	-	Supraventricular tachycardia	Supraventricular tachycardia	426761007
PVC (Ventricular premature complex)	Ventricular premature complex(es)	Ventricular premature complex	-	164884008
1AVB (First degree AV block)	-	First degree AV block	1 degree atrioventricular block	270492004
2AVB (Second degree AV block)	Second-degree AV block, Mobitz type I (Wenckebach) Second-degree AV block, Mobitz type II 2:1 AV block AV block, varying conduction AV block, advanced (high-grade)	Second degree AV block	2 degree atrioventricular block(Type one) 2 degree atrioventricular block(Type two)	54016002 28189009 164903001 195042002 284941000119107
3AVB (Third degree AV block)	AV block, complete (third-degree)	Third degree AV block	3 degree atrioventricular block	27885002
LBBB (Left bundle branch block)	Left anterior fascicular block Left posterior fascicular block Left bundle-branch block	Left anterior fascicular block Left posterior fascicular block Complete left bundle branch block	Left bundle branch block	445118002 445211001 164909002
RBBB (Right bundle branch block)	Incomplete right bundle-branch block Right bundle-branch block	Incomplete right bundle branch block Complete right bundle branch block	Right bundle branch block	713426002 59118001 164907000
LAO/LAE (Left atrial overload/enlargement)	Left atrial enlargement	Left atrial overload/enlargement	-	67741000119109
LVH (Left ventricular hypertrophy)	Left ventricular hypertrophy	Left ventricular hypertrophy	-	164873001
RVH (Right ventricular hypertrophy)	Right ventricular hypertrophy	Right ventricular hypertrophy	-	-
AMI (Anterior myocardial infarction)	Anterior MI	Anterior myocardial infarction	-	-
IMI (Inferior myocardial infarction)	Inferior MI	Inferior myocardial infarction	-	-
ASMI (Anteroseptal myocardial infarction)	Anteroseptal MI	Anteroseptal myocardial infarction	-	-

2. **PTB-XL:** <https://physionet.org/content/ptb-xl/1.0.3/>
3. **SXPH:** <https://physionet.org/content/ecg-arrhythmia/1.0.0/>
4. **G12EC:** <https://physionet.org/content/challenge-2020/1.0.2/>

### C.3.2 Preprocessing

Raw 12-lead ECG signals have varying sequence lengths and raw 12-lead ECG signals have varying sequence lengths and annotated standards which must be standardized before FL training. Therefore, we first set a signal length to 10 seconds. We pad the signal with edge value at the edge for those whose length is shorter than 10 seconds and cut off the signal at 10 seconds for those whose length is longer than 10 seconds. Next, we only save the records whose label occurs in at least two databases. Finally, we align the labels of records in different databases. The relationship between the original label and our label is shown in Table 7.

### C.4 Baseline, loss function and evaluation

**Baseline Model.** We implement a ResNet1d model with 34 layers. The final layer output is passed through a sigmoid function to encode the probability that each label corresponds to one 12-lead ECG signal.

**Loss function.** The model was directly trained for the Binary CrossEntropy Loss (BCELoss), defined as:

$$\text{BCE}(\mathbf{y}, \hat{\mathbf{y}}) = -[\sum_{i=1}^n y_i \log(\hat{y}_i) + \sum_{i=1}^n (1 - y_i) \log(1 - \hat{y}_i)] \quad (1)$$

**Evaluation Metrics.** In multi-label classification for Fed-ECG, the micro F1 score is used as the main metric to evaluate the performance of the model. Given  $N$  labels, the micro-precision ( $P_{\text{micro}}$ ) and micro-recall ( $R_{\text{micro}}$ ) are calculated as  $P_{\text{micro}} = \frac{\sum_{i=1}^N \text{TP}_i}{\sum_{i=1}^N (\text{TP}_i + \text{FP}_i)}$  and  $R_{\text{micro}} = \frac{\sum_{i=1}^N \text{TP}_i}{\sum_{i=1}^N (\text{TP}_i + \text{FN}_i)}$ , where  $\text{TP}_i$  is the number of true positives for label  $i$ ,  $\text{FP}_i$  is the number of false positives for label  $i$ ,  $\text{FN}_i$  is the number of false negatives for label  $i$ . The micro F1 score ( $F1_{\text{micro}}$ ) is then calculated as:

$$F1_{\text{micro}} = \frac{2 \cdot P_{\text{micro}} \cdot R_{\text{micro}}}{P_{\text{micro}} + R_{\text{micro}}} \quad (2)$$

For Fed-ECG's Multi-Label Classification task, the Mean Average Precision (mAP) is adopted to measure the classification performance across all labels (including long-tailed labels), calculated by averaging the average precision (AP) for each label, defined as:

$$\text{mAP} = \frac{1}{L} \sum_{i=1}^L \sum_{k=1}^n P_i(k) \Delta r_i(k) \quad (3)$$

where  $L$  is the total number of labels, and  $\text{AP}_i$  is the average precision for the  $i$ -th label,  $P_i(k)$  is the precision for label  $i$  at the  $k$ -th threshold, and  $\Delta r_i(k)$  is the change in its recall at the  $k$ -th threshold.

### C.5 Training Detail

**Optimization parameters.** We optimize the ResNet1d using SGD optimizer, with a batch size of 32. We train our model for 50 epochs on one NVIDIA A100-PCIE-40GB.

**Hyperparameter Search** For centralized and local model training, we first conduct a search for optimal learning rates from the set  $\{1e-5, 1e-4, 1e-3, 1e-2, 1e-1\}$  during centralized model training. The learning rate that yields the best micro-F1 score is then used for local model training. For the federated learning strategies, we employ the following hyperparameter grid:

- For clients’ learning rates (all strategies):  $\{1e-5, 1e-4, 1e-3, 1e-2, 1e-1\}$ .
- For server size learning rate (Scaffold strategy only):  $\{1e-2, 1e-1, 1.0\}$ .
- For FedProx and Ditto strategies, the parameter  $\mu$  is selected from  $\{1e-2, 1e-1, 1.0\}$ .
- For FedInit, the parameter  $\beta$  is chosen from  $\{1e-1, 1e-2, 1e-3\}$ .
- For FedSM, the parameters  $\gamma$  and  $\lambda$  are set to values from  $\{0, 0.1, 0.7, 0.9\}$  and  $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ , respectively.
- For FedALA, the parameters layer index,  $\eta$ , threshold, and num\_per\_loss are fixed at 1, 1.0, 0.1, and 10, respectively, while rand\_percent is selected from  $\{5, 50, 80\}$ .

Table 8: Hyperparameters used for the Fed-ECG.

Fed-ECG								
Methods	learning rate	optimizer	learning rate server	mu	beta	lambda	gamma	rand_percent
Central.	0.1	torch.optim.SGD	-	-	-	-	-	-
FedAvg	0.1	torch.optim.SGD	-	-	-	-	-	-
FedProx	0.1	torch.optim.SGD	-	0.01	-	-	-	-
Scaffold	0.1	torch.optim.SGD	1.0	-	-	-	-	-
FedInit	0.1	torch.optim.SGD	1.0	-	0.01	-	-	-
Ditto	0.1	torch.optim.SGD	-	0.01	-	-	-	-
FedSM	0.1	torch.optim.SGD	1.0	-	-	0.1	0	-
FedALA	0.1	torch.optim.SGD	1.0	-	-	-	-	80

**Non-IID partition.** For the non-IID partition, we first pool the training data from the four clients. Then, we cluster the samples into 10 categories based on the cosine similarity and order them according to the number of samples contained in each category. Next, the sorted samples are divided into 32 shards. finally, 8 random shards are distributed to one client. The label distribution of each client with the non-IID partition is shown in Figure 5.

### C.6 Supplementary Experiment Results

We provide additional evaluation metrics here. Table 9 presents an extensive array of evaluation metrics for various federated learning approaches applied to Fed-ECG. The Micro F1-Score (Mi-F1) and Hamming Loss (HL) serve as indicators of the overall performance, given their insensitivity to long-tail distributions. In contrast, the mean Average Precision score (mAP) provides insight into the average performance across individual labels. In addition, Figure 6 presents the evaluation metrics for each label, encompassing F1 score, precision, and recall, which more clearly demonstrates the impact of the long-tail distribution on each label.

## D Fed-ECHO

### D.1 Description

Fed-ECHO consists of three datasets: CAMUS, ECHONET-DYNAMIC, and HMC-QU. The overview of Fed-ECHO is shown in Table 5.

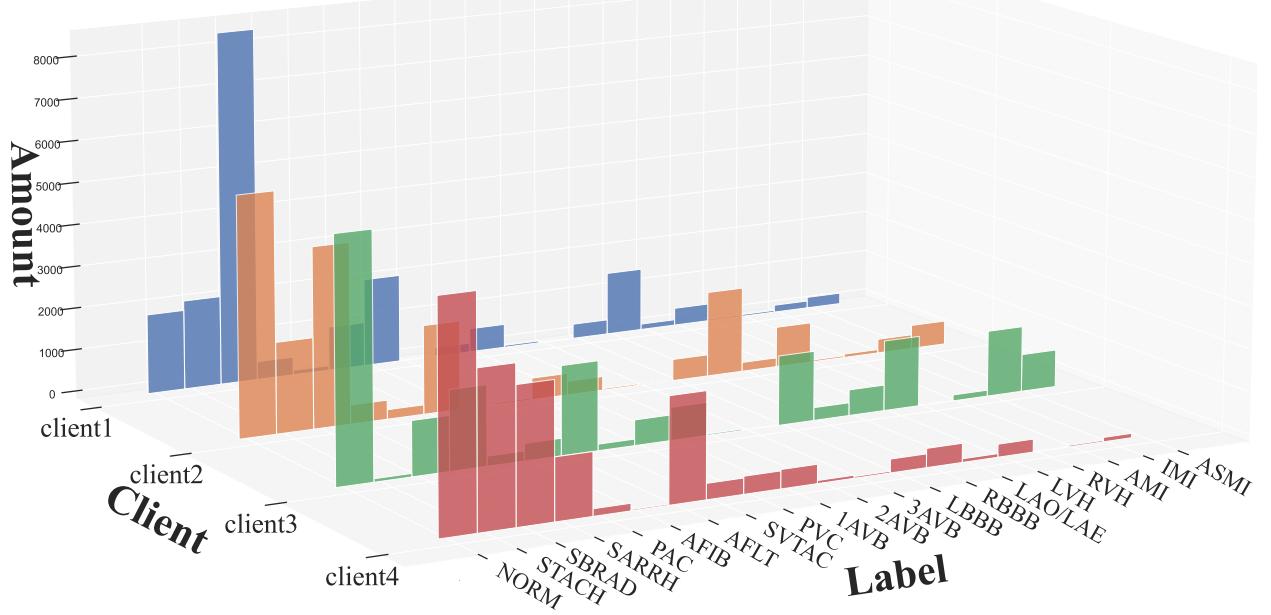


Figure 5: Label non-IID of the Fed-ECG dataset with the artificially non-IID partition, shown as the variation in the number of each label (right axis) across different clients (left axis).

Table 9: The performance of different FL methods on Fed-ECG, with Mi-F1, mAP, and HL representing Micro F1-Score, mean Average Precision score, and Hamming Loss, respectively. All metrics are present in percentage (%). The best results for each configuration are highlighted in **bold**, while the second-best results are underlined.

Methods	LOCAL												GLOBAL		
	Client1			Client2			Client3			Client4			Mi-F1↑	mAP↑	HL↓
Client1	85.8 ±1.9	58.1 ±2.6	1.5 ±0.2	52.7 ±3.4	37.8 ±2.2	5.8 ±0.4	61.5 ±1.2	19.8 ±1.2	4.4 ±0.1	49.8 ±4.2	26.7 ±3.0	6.4 ±0.6	64.3 ±2.1	32.3 ±2.0	4.1 ±0.2
Client2	69.9 ±50.0	38.9 ±30.0	3.2 ±0.1	76.8 ±90.0	55.7 ±50.0	3.1 ±0.1	26.3 ±80.0	22.7 ±30.0	9.0 ±0.2	42.2 ±80.0	31.6 ±60.0	8.1 ±0.1	50.4 ±30.0	35.9 ±70.0	6.1 ±0.1
Client3	22.7 ±0.2	29.8 ±0.7	8.2 ±0.0	17.0 ±0.4	27.2 ±0.3	10.3 ±0.2	88.1 ±0.2	37.7 ±0.4	1.3 ±0.0	56.9 ±0.4	29.4 ±0.6	5.4 ±0.1	51.5 ±0.2	32.7 ±0.2	5.5 ±0.0
Client4	23.7 ±2.0	31.7 ±2.7	8.4 ±0.9	24.7 ±3.3	30.5 ±1.5	10.1 ±1.2	61.6 ±5.5	25.3 ±2.1	5.0 ±1.2	72.3 ±10.2	38.5 ±2.8	4.1 ±1.8	44.7 ±4.3	29.3 ±2.5	7.0 ±1.1
FedAvg	69.0 ±10.1	58.5 ±1.2	3.4 ±1.1	50.3 ±5.3	54.4 ±0.5	6.2 ±0.7	77.6 ±0.7	37.2 ±0.3	2.5 ±0.1	66.3 ±0.9	39.5 ±0.5	4.2 ±0.1	67.9 ±3.8	50.8 ±0.4	3.7 ±0.5
FedProx	74.0 ±7.5	60.3 ±2.9	2.9 ±1.0	55.6 ±2.7	56.4 ±0.6	5.5 ±0.5	73.2 ±1.0	36.0 ±0.8	3.0 ±0.1	70.2 ±2.3	<b>43.8</b> ±1.8	3.8 ±0.3	68.8 ±2.6	<b>52.3</b> ±0.9	3.6 ±0.4
Scaffold	77.5 ±2.6	58.0 ±1.2	2.3 ±0.2	56.9 ±1.7	55.9 ±0.7	5.2 ±0.2	73.3 ±1.0	36.2 ±0.6	3.0 ±0.1	70.7 ±2.9	42.7 ±1.1	3.7 ±0.3	<b>70.1</b> ±0.8	<b>52.1</b> ±0.7	<b>3.4</b> ±0.1
FedInit	73.0 ±6.6	58.2 ±0.7	3.1 ±1.0	54.1 ±5.2	55.6 ±1.3	5.9 ±1.3	73.5 ±0.9	36.6 ±0.5	3.0 ±0.1	67.8 ±2.9	41.5 ±1.1	4.1 ±0.3	68.1 ±0.3	51.5 ±3.0	3.8 ±0.5
Ditto	82.8 ±4.4	<b>63.1</b> ±4.2	1.8 ±0.4	<b>74.8</b> ±1.4	<b>58.3</b> ±0.6	<b>3.5</b> ±0.2	<b>86.5</b> ±1.5	<b>38.1</b> ±0.6	<b>1.5</b> ±0.2	<b>73.4</b> ±6.7	42.2 ±4.0	<b>3.6</b> ±0.9	68.1 ±2.9	48.7 ±1.4	<b>3.6</b> ±0.3
FedSM	77.2 ±7.2	58.8 ±1.3	2.3 ±0.6	59.1 ±4.5	56.4 ±1.4	5.1 ±0.5	69.8 ±0.8	35.0 ±0.5	3.5 ±0.1	67.7 ±3.6	42.9 ±2.4	4.1 ±0.4	68.9 ±2.5	51.2 ±0.7	3.6 ±0.3
FedALA	<b>84.4</b> ±4.0	62.0 ±7.0	<b>1.6</b> ±0.4	71.7 ±5.7	<b>57.1</b> ±2.2	<b>3.8</b> ±0.6	<b>88.2</b> ±0.1	<b>37.4</b> ±0.2	<b>1.3</b> ±0.0	66.7 ±5.9	41.2 ±2.3	4.4 ±0.7	<b>67.8</b> ±1.9	50.8 ±1.3	3.7 ±0.3
Central.	84.9 ±0.5	54.8 ±0.5	1.6 ±0.1	71.4 ±5.0	55.2 ±2.9	3.8 ±0.6	84.1 ±1.6	36.5 ±1.1	1.7 ±0.2	72.2 ±3.7	41.5 ±1.3	3.6 ±0.3	80.0 ±2.1	63.2 ±2.8	2.3 ±0.2

**CAMUS.** This database consists of clinical exams from 500 patients, acquired at the University Hospital of St Etienne (France). All images are labeled with three areas: endocardium of the left ventricle ( $LV_{Endo}$ ), epicardium of the left ventricle ( $LV_{Epi}$ ), and left atrium wall (LA). The image size varies from  $584 \times 354$  to  $1945 \times 1181$ .

**ECHONET-DYNAMIC.** This database contains 10,0230 echocardiogram videos where two frames are annotated with only  $LV_{Endo}$  area. All frames are resized to  $112 \times 112$ .

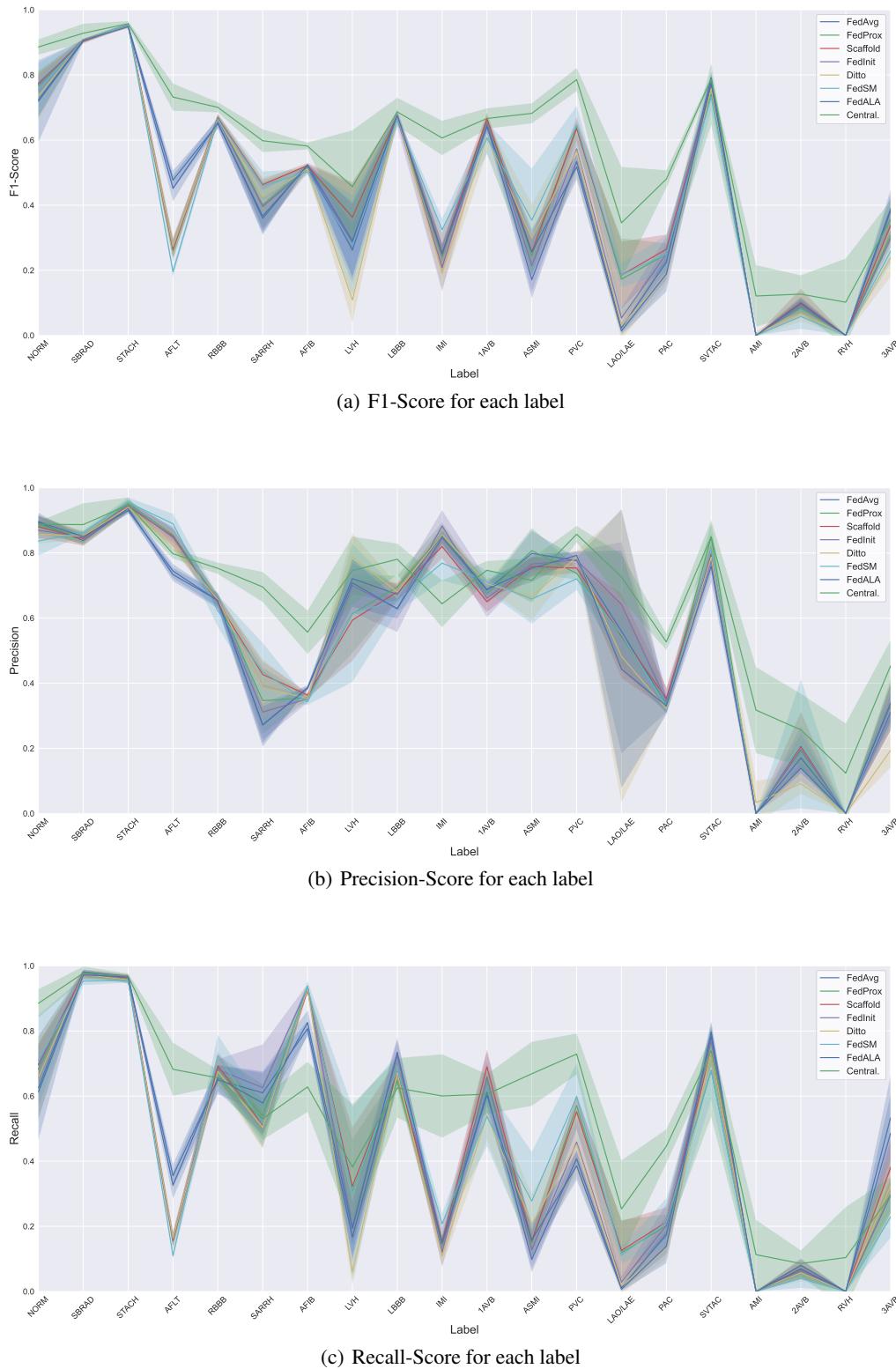


Figure 6: Evaluation metrics for each label on Fed-ECG among different FL methods.

**HMC-QU.** This database contains 109 echocardiogram videos collected at the Hamad Medical Corporation Hospital in Qatar. The frames of one cardiac cycle in each video are annotated with  $LV_{Epi}$  area. The video frame size varies from  $422 \times 636$  to  $768 \times 1024$  while all labels are resized to  $224 \times 224$ .

## D.2 License and Ethics

Both CAMUS and HMC-QU datasets are open-access. HMC-QU database requires the user to have a Kaggle account, while the ECHONET-DYNAMIC database requires the user to have a Stanford AIMI account and to accept its agreement. It is licensed under the Stanford University Dataset Research Use Agreement.

## D.3 Download and preprocessing

### D.3.1 Download

The three datasets can be downloaded using the URLs below:

1. **CAMUS:** <https://humanheart-project.creatis.insa-lyon.fr/database/#collection/6373703d73e9f0047faa1bc8>
2. **ECHONET-DYNAMIC:** <https://echonet.github.io/dynamic/index.html#access>
3. **HMC-QU:** <https://www.kaggle.com/datasets/aysendegerli/hmcqu-dataset/data>

### D.3.2 Preprocessing

Raw echocardiograms have varying frame sizes, modalities, and mask labels, which must be standardized before training. Therefore, as a first step, we extract frames that are annotated and store them as images. We then resize them to a common ( $112 \times 112$ ) shape. Finally, we align the labels of records in different databases. We use 1, 2, 3 representing  $LV_{Endo}$ ,  $LV_{Epi}$  and LA respectively. The samples of Fed-ECHO are shown in Figure 7.



(a) Sample from Institution 1. (b) Sample from Institution 2. (c) ample from Institution 3.

Figure 7: Echocardiogram of each institution in Fed-ECHO.  $LV_{Endo}$ ,  $LV_{Epi}$  and LA are shown in red, green and blue respectively.

## D.4 Baseline, loss function and evaluation

**Baseline Model.** A U-net architecture is employed in this study, utilizing echocardiographic images as input to forecast masks delineating four distinct cardiac regions. The U-net model represents a conventional convolutional neural network design frequently deployed in the realm of biomedical image segmentation endeavors. Its application is tailored towards semantic segmentation, a process wherein individual pixels within an image are categorized based on semantic content.

**Loss function.** We use a CrossEntropy Loss (CELoss) for training. Note that, for centralized supervised learning and client training in FedAvg, FedProx, Scaffold, and Ditto strategies, we ignore label with value 0 when calculating CELoss for data from client 2 or 3, since region with label 0 may not be true ground truth in these clients.

**Evaluation Metrics.** We use the Dice similarity index and 2D Hausdorff distance ( $d_H$ ) to measure the accuracy of the segmentation output. Dice index is calculated as:

$$\text{DICE}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{2 \sum_{i=1}^n y_i \hat{y}_i}{\sum_{i=1}^n y_i + \sum_{i=1}^n \hat{y}_i} \quad (4)$$

The Hausdorff distance is calculated as:

$$d_H(\mathbf{y}, \hat{\mathbf{y}}) = \max\{d(\mathbf{y}, \hat{\mathbf{y}}), d(\hat{\mathbf{y}}, \mathbf{y})\}, \quad (5)$$

where  $d(\mathbf{y}, \hat{\mathbf{y}})$  represents the minimum distance among points at the edge of  $\mathbf{y}$  and points at the edge of  $\hat{\mathbf{y}}$ .

Note that, to better measure the model segmentation performance, for clients 2, and 3, we select only 200 labeled frames for testing.

## D.5 Training Detail

**Optimization parameters.** We optimize our model using the SGD optimizer, with a batch size of 32. We train our model for 50 epochs on one NVIDIA A100-PCIE-40GB.

**Hyperparameter Search** For centralized and local model training, we first explore learning rates from the set {1e-4, 1e-3, 1e-2, 1e-1.5, 1e-1} during centralized model training. The learning rate that achieves the best Dice index is then utilized for local model training. For the federated learning strategies, we employ the following hyperparameter grid:

- For clients' learning rates (all strategies except Fed-Consist): {1e-4, 1e-3, 1e-2, 1e-1.5, 1e-1}.
- For server size learning rate (Scaffold strategy only): {1e-2, 1e-1, 1.0}.
- For FedProx and Ditto strategies, the parameter  $\mu$  is selected from {1e-2, 1e-1, 1.0}.
- For FedInit, the parameter  $\beta$  is chosen from {1e-1, 1e-2, 1e-3}.
- For FedSM, the parameters  $\gamma$  and  $\lambda$  are set to {0, 0.1, 0.7, 0.9} and {0.1, 0.3, 0.5, 0.7, 0.9}, respectively.
- For FedALA, the parameters layer index,  $\eta$ , threshold, and num\_per\_loss are fixed at 1, 1.0, 0.1, and 10, respectively, while rand\_percent is chosen from {5, 50, 80}.

For Fed-Consist, we introduce Gaussian noise with a variance of 0.1 as augmentation. The learning rates for labeled clients are searched from {1e-2, 1e-3, 1e-4}, while those for unlabeled clients are explored within {1e-3, 1e-4, 1e-5, 5e-6, 1e-6}. The parameter  $\tau$  is varied from {0.5, 0.7, 0.9}.

Additionally, for FedPSL, we further search the parameters  $\alpha$  and  $\beta$  from {1e-0.5, 1e-1, 1e-1.5, 1e-2, 1e-3} and {1e-1, 1e-1.5, 1e-2, 1e-3, 1e-4, 1e-5}, respectively. The optimal values found are  $\alpha = 1e - 1.5$  and  $\beta = 1e - 5$ .

Table 10: Hyperparameters used for the Fed-ECHO.

Fed-ECHO									
Methods	learning rate	optimizer	learning rate server	mu	beta	lambda	gamma	rand_percent	$\tau$
Central.(sup)	0.1	torch.optim.SGD	-	-	-	-	-	-	-
Central.(ssup)	0.1	torch.optim.SGD	-	-	-	-	-	-	-
FedAvg	0.1	torch.optim.SGD	-	-	-	-	-	-	-
FedProx	0.1	torch.optim.SGD	-	0.1	-	-	-	-	-
Scaffold	0.1	torch.optim.SGD	1.0	-	-	-	-	-	-
FedInit	0.1	torch.optim.SGD	1.0	-	1e-2	-	-	-	-
Ditto	0.1	torch.optim.SGD	-	0.1	-	-	-	-	-
FedSM	0.1	torch.optim.SGD	1.0	-	-	0.1	0	-	-
FedALA	0.1	torch.optim.SGD	1.0	-	-	-	-	5	-
FedPSL	0.1	torch.optim.SGD	1.0	-	1e-5	-	-	-	-
Fed-Consist	0.0001(labeled client) 1e-6(unlabeled client)	torch.optim.SGD	-	-	-	-	-	-	0.9